

# Introduction to Convolutional Neural Networks

Jianxin Wu

LAMDA Group

National Key Lab for Novel Software Technology

Nanjing University, China

wujx2001@gmail.com

March 21, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Preliminaries</b>	<b>3</b>
2.1	Tensor and vectorization . . . . .	3
2.2	Vector calculus and the chain rule . . . . .	4
<b>3</b>	<b>CNN in a nutshell</b>	<b>5</b>
3.1	The architecture . . . . .	5
3.2	The forward run . . . . .	6
3.3	Stochastic gradient descent (SGD) . . . . .	6
3.4	Error back propagation . . . . .	8
<b>4</b>	<b>Layer input, output and notations</b>	<b>9</b>
<b>5</b>	<b>The ReLU layer</b>	<b>10</b>
<b>6</b>	<b>The convolution layer</b>	<b>11</b>
6.1	What is convolution? . . . . .	11
6.2	Why to convolve? . . . . .	13
6.3	Convolution as matrix product . . . . .	15
6.4	The Kronecker product . . . . .	17
6.5	Backward propagation: update the parameters . . . . .	17
6.6	Even higher dimensional indicator matrices . . . . .	19
6.7	Backward propagation: prepare supervision signal for the previous layer . . . . .	20
6.8	Fully connected layer as a convolution layer . . . . .	22
<b>7</b>	<b>The pooling layer</b>	<b>23</b>

<b>8 A case study: the VGG-16 net</b>	<b>25</b>
8.1 VGG-Verydeep-16 . . . . .	25
8.2 Receptive field . . . . .	27
<b>9 Remarks</b>	<b>28</b>
<b>Exercises</b>	<b>28</b>

# 1 Introduction

This is a note that describes how a Convolutional Neural Network (CNN) operates from a mathematical perspective. This note is self-contained, and the focus is to make it comprehensible to beginners in the CNN field.

The Convolutional Neural Network (CNN) has shown excellent performance in many computer vision and machine learning problems. Many solid papers have been published on this topic, and quite some high quality open source CNN software packages have been made available.

There are also well-written CNN tutorials or CNN software manuals. However, I believe that an introductory CNN material specifically prepared for beginners is still needed. Research papers are usually very terse and lack details. It might be difficult for beginners to read such papers. A tutorial targeting experienced researchers may not cover all the necessary details to understand how a CNN runs.

This note tries to present a document that

- is self-contained. It is expected that all required mathematical background knowledge are introduced in this note itself (or in other notes for this course);
- has details for all the derivations. This note tries to explain all the necessary math in details. We try not to ignore an important step in a derivation. Thus, it should be possible for a beginner to follow (although an expert may feel this note tautological.)
- ignores implementation details. The purpose is for a reader to understand how a CNN runs at the mathematical level. We will ignore those implementation details. In CNN, making correct choices for various implementation details is one of the keys to its high accuracy (that is, “the devil is in the details”). However, we intentionally left this part out, in order for the reader to focus on the mathematics. After understanding the mathematical principles and details, it is more advantageous to learn these implementation and design details with hands-on experience by playing with CNN programming.

CNN is useful in a lot of applications, especially in image related tasks. Applications of CNN include image classification, image semantic segmentation,

object detection in images, etc. We will focus on image classification (or categorization) in this note. In image categorization, every image has a major object which occupies a large portion of the image. An image is classified into one of the classes based on the identity of its main object, e.g., dog, airplane, bird, etc.

## 2 Preliminaries

We start by a discussion of some background knowledge that are necessary in order to understand how a CNN runs. One can ignore this section if he/she is familiar with these basics.

### 2.1 Tensor and vectorization

Everybody is familiar with vectors and matrices. We use a symbol shown in boldface to represent a vector, e.g.,  $\mathbf{x} \in \mathbb{R}^D$  is a column vector with  $D$  elements. We use a capital letter to denote a matrix, e.g.,  $X \in \mathbb{R}^{H \times W}$  is a matrix with  $H$  rows and  $W$  columns. The vector  $\mathbf{x}$  can also be viewed as a matrix with 1 column and  $D$  rows.

These concepts can be generalized to higher-order matrices, i.e., tensors. For example,  $\mathbf{x} \in \mathbb{R}^{H \times W \times D}$  is an order 3 (or third order) tensor. It contains  $HW D$  elements, and each of them can be indexed by an index triplet  $(i, j, d)$ , with  $0 \leq i < H$ ,  $0 \leq j < W$ , and  $0 \leq d < D$ . Another way to view an order 3 tensor is to treat it as containing  $D$  channels of matrices. Every channel is a matrix with size  $H \times W$ . The first channel contains all the numbers in the tensor that are indexed by  $(i, j, 0)$ . When  $D = 1$ , an order 3 tensor reduces to a matrix.

We have interacted with tensors day-to-day. A scalar value is a zeroth-order (order 0) tensor; a vector is an order 1 tensor; and a matrix is a second order tensor. A color image is in fact an order 3 tensor. An image with  $H$  rows and  $W$  columns is a tensor with size  $H \times W \times 3$ : if a color image is stored in the RGB format, it has 3 channels (for R, G and B, respectively), and each channel is a  $H \times W$  matrix (second order tensor) that contains the R (or G, or B) values of all pixels.

It is beneficial to represent images (or other types of raw data) as a tensor. In early computer vision and pattern recognition, a color image (which is an order 3 tensor) is often converted to the gray-scale version (which is a matrix) because we know how to handle matrices much better than tensors. The color information is lost during this conversion. But color is very important in various image (or video) based learning and recognition problems, and we do want to process color information in a principled way, e.g., as in CNN.

Tensors are essential in CNN. **The input, intermediate representation, and parameters in a CNN are all tensors.** Tensors with order higher than 3 are also widely used in a CNN. For example, we will soon see that the convolution kernels in a convolution layer of a CNN form an order 4 tensor.

Given a tensor, we can arrange all the numbers inside it into a long vector, following a pre-specified order. For example, in Matlab, the  $(:)$  operator

converts a matrix into a column vector in the column-first order. An example is:

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad A(\cdot) = (1, 3, 2, 4)^T = \begin{bmatrix} 1 \\ 3 \\ 2 \\ 4 \end{bmatrix}. \quad (1)$$

In mathematics, we use the notation “vec” to represent this vectorization operator. That is,  $\text{vec}(A) = (1, 3, 2, 4)^T$  in the example in Equation 1. In order to vectorize an order 3 tensor, we could vectorize its first channel (which is a matrix and we already know how to vectorize it), then the second channel, ..., till all channels are vectorized. The vectorization of the order 3 tensor is then the concatenation of the vectorization of all the channels in this order.

The vectorization of an order 3 tensor is a recursive process, which utilizes the vectorization of order 2 tensors. This recursive process can be applied to vectorize an order 4 (or even higher order) tensor in the same manner.

## 2.2 Vector calculus and the chain rule

The CNN learning process depends on vector calculus and the chain rule. Suppose  $z$  is a scalar (i.e.,  $z \in \mathbb{R}$ ) and  $\mathbf{y} \in \mathbb{R}^H$  is a vector. If  $z$  is a function of  $\mathbf{y}$ , then the partial derivative of  $z$  with respect to  $\mathbf{y}$  is a vector, defined as

$$\left[ \frac{\partial z}{\partial \mathbf{y}} \right]_i = \frac{\partial z}{\partial y_i}. \quad (2)$$

In other words,  $\frac{\partial z}{\partial \mathbf{y}}$  is a vector having *the same size* as  $\mathbf{y}$ , and its  $i$ -th element is  $\frac{\partial z}{\partial y_i}$ . Also note that  $\frac{\partial z}{\partial \mathbf{y}^T} = \left( \frac{\partial z}{\partial \mathbf{y}} \right)^T$ .

Furthermore, suppose  $\mathbf{x} \in \mathbb{R}^W$  is another vector, and  $\mathbf{y}$  is a function of  $\mathbf{x}$ . Then, the partial derivative of  $\mathbf{y}$  with respect to  $\mathbf{x}$  is defined as

$$\left[ \frac{\partial \mathbf{y}}{\partial \mathbf{x}^T} \right]_{ij} = \frac{\partial y_i}{\partial x_j}. \quad (3)$$

This **partial derivative is a  $H \times W$  matrix**, whose entry at the intersection of the  $i$ -th row and  $j$ -th column is  $\frac{\partial y_i}{\partial x_j}$ .

It is easy to see that  $z$  is a function of  $\mathbf{x}$  in a chain-like argument: a function maps  $\mathbf{x}$  to  $\mathbf{y}$ , and another function maps  $\mathbf{y}$  to  $z$ . The chain rule can be used to compute  $\frac{\partial z}{\partial \mathbf{x}^T}$ , as

$$\frac{\partial z}{\partial \mathbf{x}^T} = \frac{\partial z}{\partial \mathbf{y}^T} \frac{\partial \mathbf{y}}{\partial \mathbf{x}^T}. \quad (4)$$

A sanity check for Equation 4 is to check the matrix / vector dimensions. Note that  $\frac{\partial z}{\partial \mathbf{y}^T}$  is a row vector with  $H$  elements, or a  $1 \times H$  matrix. (Be reminded that  $\frac{\partial z}{\partial \mathbf{y}}$  is a column vector). Since  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}^T}$  is an  $H \times W$  matrix, the vector / matrix multiplication between them is valid, and the result should be a row vector with  $W$  elements, which matches the dimensionality of  $\frac{\partial z}{\partial \mathbf{x}^T}$ .

For specific rules to calculate partial derivatives of vectors and matrices, please refer to the Matrix Cookbook.

### 3 CNN in a nutshell

In this section, we will see how a CNN trains and predicts in the abstract level, with the details left out for later sections.

#### 3.1 The architecture

A CNN usually takes an order 3 tensor as its input, e.g., an image with  $H$  rows,  $W$  columns, and 3 channels (R, G, B color channels). Higher order tensor inputs, however, can be handled by CNN in a similar fashion. The input then sequentially goes through a series of processing. One processing step is usually called a layer, which could be a convolution layer, a pooling layer, a normalization layer, a fully connected layer, a loss layer, etc. We will introduce the details of these layers later in this note.<sup>1</sup>

For now, let us give an abstract description of the CNN structure first.

$$\mathbf{x}^1 \longrightarrow \boxed{\mathbf{w}^1} \longrightarrow \mathbf{x}^2 \longrightarrow \dots \longrightarrow \mathbf{x}^{L-1} \longrightarrow \boxed{\mathbf{w}^{L-1}} \longrightarrow \mathbf{x}^L \longrightarrow \boxed{\mathbf{w}^L} \longrightarrow z \quad (5)$$

The above Equation 5 illustrates how a CNN runs layer by layer in a forward pass. The input is  $\mathbf{x}^1$ , usually an image (order 3 tensor). It goes through the processing in the first layer, which is the first box. We denote the parameters involved in the first layer's processing collectively as a tensor  $\mathbf{w}^1$ . The output of the first layer is  $\mathbf{x}^2$ , which also acts as the input to the second layer processing.

This processing proceeds till all layers in the CNN has been finished, which outputs  $\mathbf{x}^L$ . One additional layer, however, is added for backward error propagation, a method that learns good parameter values in the CNN. Let's suppose the problem at hand is an image classification problem with  $C$  classes. A commonly used strategy is to output  $\mathbf{x}^L$  as a  $C$  dimensional vector, whose  $i$ -th entry encodes the prediction (posterior probability of  $\mathbf{x}^1$  comes from the  $i$ -th class). To make  $\mathbf{x}^L$  a probability mass function, we can set the processing in the  $(L-1)$ -th layer as a softmax transformation of  $\mathbf{x}^{L-1}$  (cf. the distance metric and data transformation note). In other applications, the output  $\mathbf{x}^L$  may have other forms and interpretations.

The last layer is a loss layer. Let us suppose  $\mathbf{t}$  is the corresponding target (ground-truth) value for the input  $\mathbf{x}^1$ , then a cost or loss function can be used to measure the discrepancy between the CNN prediction  $\mathbf{x}^L$  and the target  $\mathbf{t}$ . For example, a simple loss function could be

$$z = \frac{1}{2} \|\mathbf{t} - \mathbf{x}^L\|^2, \quad (6)$$

---

<sup>1</sup>We will give detailed introductions to three types of layers: convolution, pooling, and ReLU, which are the key parts of almost all CNN models. Proper normalization, e.g., batch normalization or cross-layer normalization is important in the optimization process for learning good parameters in a CNN. I may add these contents in the next update.

although more complex loss functions are usually used. This squared  $\ell_2$  loss can be used in a regression problem. In a classification problem, **the cross entropy loss** is often used. The ground-truth in a classification problem is a categorical variable  $t$ . We first convert the categorical variable  $t$  to a  $C$  dimensional vector  $\mathbf{t}$  (cf. the distance metric and data transformation note). Now both  $\mathbf{t}$  and  $\mathbf{x}^L$  are probability mass functions, and the cross entropy loss measures the distance between them. Hence, we can minimize the cross entropy (cf. the information theory note.) Equation 5 explicitly models the loss function as a loss layer, whose processing is modeled as a box with parameters  $\mathbf{w}^L$ .

Note that some layers may not have any parameters, that is,  $\mathbf{w}^i$  may be empty for some  $i$ . The softmax layer is one such example.

### 3.2 The forward run

Suppose all the parameters of a CNN model  $\mathbf{w}^1, \dots, \mathbf{w}^{L-1}$  have been learned, then we are ready to use this model for prediction. Prediction only involves running the CNN model forward, i.e., in the direction of the arrows in Equation 5.

Let's take the image classification problem as an example. Starting from the input  $\mathbf{x}^1$ , we make it pass the processing of the first layer (the box with parameters  $\mathbf{w}^1$ ), and get  $\mathbf{x}^2$ . In turn,  $\mathbf{x}^2$  is passed into the second layer, etc. Finally, we achieve  $\mathbf{x}^L \in \mathbb{R}^C$ , which estimates the posterior probabilities of  $\mathbf{x}^1$  belonging to the  $C$  categories. We can output the CNN prediction as

$$\arg \max_i x_i^L. \quad (7)$$

Note that the loss layer is not needed in prediction. It is only useful when we try to learn CNN parameters using a set of training examples. Now, the problem is: how do we learn the model parameters?

### 3.3 Stochastic gradient descent (SGD)

As in many other learning systems, the parameters of a CNN model are optimized to minimize **the loss  $z$** , i.e., we want the prediction of a CNN model to match the ground-truth labels.

Let's suppose one training example  $\mathbf{x}^1$  is given for training such parameters. The training process involves running the CNN network in both directions. We first run the network in the forward pass to get  $\mathbf{x}^L$  to achieve a prediction using the current CNN parameters. Instead of outputting a prediction, we need to compare the prediction with the target  $\mathbf{t}$  corresponding to  $\mathbf{x}^1$ , that is, continue running the forward pass till the last loss layer. Finally, we achieve a loss  $z$ .

The loss  $z$  is then a supervision signal, guiding how the parameters of the model should be modified (updated). And the SGD way of modifying the parameters is

$$\mathbf{w}^i \longleftarrow \mathbf{w}^i - \eta \frac{\partial z}{\partial \mathbf{w}^i}. \quad (8)$$

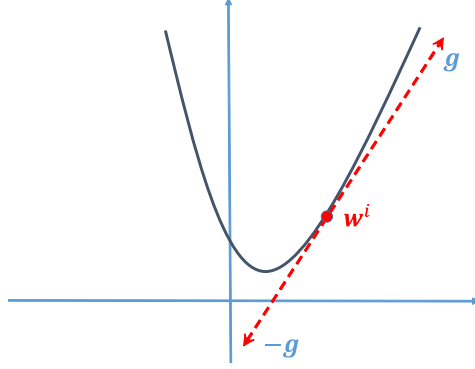


Figure 1: Illustration of the gradient descent method.

*A cautious note about the notation.* In most CNN materials, a superscript indicates the “time” (e.g., training epochs). But in this note, we use the superscript to denote the layer index. Please do not get confused. We do not use an additional index variable to represent time. In Equation 8, the  $\leftarrow$  sign implicitly indicates that the parameters  $\mathbf{w}^i$  (of the  $i$ -layer) are updated from time  $t$  to  $t + 1$ . If a time index  $t$  is explicitly used, this equation will look like

$$(\mathbf{w}^i)^{t+1} = (\mathbf{w}^i)^t - \eta \frac{\partial z}{\partial (\mathbf{w}^i)^t}. \quad (9)$$

In Equation 8, the partial derivative  $\frac{\partial z}{\partial \mathbf{w}^i}$  measures the rate of increase of  $z$  with respect to the changes in different dimensions of  $\mathbf{w}^i$ . This partial derivative vector is called the *gradient* in mathematical optimization. Hence, in a small local region around the current value of  $\mathbf{w}^i$ , to move  $\mathbf{w}^i$  in the direction determined by the gradient will increase the objective value  $z$ . In order to minimize the loss function, we should update  $\mathbf{w}^i$  along the opposite direction of the gradient. This updating rule is called the gradient descent. Gradient descent is illustrated in Figure 1, in which the gradient is denoted by  $\mathbf{g}$ .

If we move too far in the negative gradient direction, however, the loss function may increase. Hence, in every update we only change the parameters by a small proportion of the negative gradient, controlled by  $\eta$  (the learning rate).  $\eta > 0$  is usually set to a small number (e.g.,  $\eta = 0.001$ ). One update based on  $\mathbf{x}^1$  will make the loss smaller for this particular training example if the learning rate is not too large. However, it is very possible that it will make the loss of some other training examples become larger. Hence, we need to update the parameters using all training examples. When all training examples have been used to update the parameters, we say one *epoch* has been processed. One epoch will in general reduce the average loss on the training set until the learning system overfits the training data. Hence, we can repeat the gradient descent updating epochs and terminate at some point to obtain the CNN parameters (e.g., we can terminate when the average loss on a validation set increases).

Gradient descent may seem simple in its math form (Equation 8), but it is a very tricky operation in practice. For example, if we update the parameters using only gradient calculated from only one training example, we will observe an unstable loss function: the average loss of all training examples will bounce up and down at very high frequency. This is because the gradient is estimated using only one training example instead of the entire training set. Updating the parameters using the gradient estimated from a (usually) small subset of training examples is called the *stochastic gradient descent*. Contrary to single example based SGD, we can compute the gradient using all training examples and then update the parameters. However, this *batch* processing requires a lot of computations because the parameters are updated only once in an epoch, and is hence impractical, especially when the number of training examples is large.

A compromise is to use a *mini-batch* of training examples, to compute the gradient using this mini-batch, and to update the parameters correspondingly. For example, we can set 32 or 64 examples as a mini-batch. Stochastic gradient descent (SGD) (using the mini-batch strategy) is the mainstream method to learn a CNN’s parameters. We also want to note that when mini-batch is used, the input of the CNN becomes an order 4 tensor, e.g.,  $H \times W \times 3 \times 32$  if the mini-batch size is 32.

A new problem now becomes apparent: how to compute the gradient, which seems a very complex task?

### 3.4 Error back propagation

The last layer’s partial derivatives are easy to compute. Because  $\mathbf{x}^L$  is connected to  $z$  directly under the control of parameters  $\mathbf{w}^L$ , it is easy to compute  $\frac{\partial z}{\partial \mathbf{w}^L}$ . This step is only needed when  $\mathbf{w}^L$  is not empty. In the same spirit, it is also easy to compute  $\frac{\partial z}{\partial \mathbf{x}^L}$ . For example, if the squared  $\ell_2$  loss is used, we have an empty  $\frac{\partial z}{\partial \mathbf{w}^L}$ , and  $\frac{\partial z}{\partial \mathbf{x}^L} = \mathbf{x}^L - \mathbf{t}$ .

In fact, for every layer, we compute two sets of gradients: the partial derivatives of  $z$  with respect to the layer parameters  $\mathbf{w}^i$ , and that layer’s input  $\mathbf{x}^i$ .

- The term  $\frac{\partial z}{\partial \mathbf{w}^i}$ , as seen in Equation 8, can be used to update the current ( $i$ -th) layer’s parameters;
- The term  $\frac{\partial z}{\partial \mathbf{x}^i}$  can be used to update parameters backwards, e.g., to the  $(i - 1)$ -th layer. An intuitive explanation is:  $\mathbf{x}^i$  is the output of the  $(i - 1)$ -th layer and  $\frac{\partial z}{\partial \mathbf{x}^i}$  is how  $\mathbf{x}^i$  should be changed to reduce the loss function. Hence, we could view  $\frac{\partial z}{\partial \mathbf{x}^i}$  as the part of the “error” supervision information propagated from  $z$  backward till the current layer, in a layer by layer fashion. Thus, we can continue the back propagation process, and use  $\frac{\partial z}{\partial \mathbf{x}^i}$  to propagate the errors backward to the  $(i - 1)$ -th layer.

This layer-by-layer backward updating procedure makes learning a CNN much easier.

Let’s take the  $i$ -th layer as an example. When we are updating the  $i$ -th layer, the back propagation process for the  $(i + 1)$ -th layer must have been finished.



That is, we already computed the terms  $\frac{\partial z}{\partial \mathbf{w}^{i+1}}$  and  $\frac{\partial z}{\partial \mathbf{x}^{i+1}}$ . Both are stored in memory and ready for use.

Now our task is to compute  $\frac{\partial z}{\partial \mathbf{w}^i}$  and  $\frac{\partial z}{\partial \mathbf{x}^i}$ . Using the chain rule, we have

$$\frac{\partial z}{\partial (\text{vec}(\mathbf{w}^i)^T)} = \frac{\partial z}{\partial (\text{vec}(\mathbf{x}^{i+1})^T)} \frac{\partial \text{vec}(\mathbf{x}^{i+1})}{\partial (\text{vec}(\mathbf{w}^i)^T)}, \quad (10)$$

$$\frac{\partial z}{\partial (\text{vec}(\mathbf{x}^i)^T)} = \frac{\partial z}{\partial (\text{vec}(\mathbf{x}^{i+1})^T)} \frac{\partial \text{vec}(\mathbf{x}^{i+1})}{\partial (\text{vec}(\mathbf{x}^i)^T)}. \quad (11)$$

Since  $\frac{\partial z}{\partial \mathbf{x}^{i+1}}$  is already computed and stored in memory, it requires just a matrix reshaping operation (vec) and an additional transpose operation to get  $\frac{\partial z}{\partial (\text{vec}(\mathbf{x}^{i+1})^T)}$ , which is the first term in the right hand side (RHS) of both equations. So long as we can compute  $\frac{\partial \text{vec}(\mathbf{x}^{i+1})}{\partial (\text{vec}(\mathbf{w}^i)^T)}$  and  $\frac{\partial \text{vec}(\mathbf{x}^{i+1})}{\partial (\text{vec}(\mathbf{x}^i)^T)}$ , we can easily get what we want (the left hand side of both equations).

$\frac{\partial \text{vec}(\mathbf{x}^{i+1})}{\partial (\text{vec}(\mathbf{w}^i)^T)}$  and  $\frac{\partial \text{vec}(\mathbf{x}^{i+1})}{\partial (\text{vec}(\mathbf{x}^i)^T)}$  are much easier to compute than directly computing  $\frac{\partial z}{\partial (\text{vec}(\mathbf{w}^i)^T)}$  and  $\frac{\partial z}{\partial (\text{vec}(\mathbf{x}^i)^T)}$ , because  $\mathbf{x}^i$  is directly related to  $\mathbf{x}^{i+1}$ , through a function with parameters  $\mathbf{w}^i$ . The details of these partial derivatives will be discussed in the following sections.

## 4 Layer input, output and notations

Now that the CNN architecture is clear, we will discuss in detail the different types of layers, starting from the ReLU layer, which is the simplest layer among those we discuss in this note. But before we start, we need to further refine our notations.

Suppose we are considering the  $l$ -th layer, whose inputs form an order 3 tensor  $\mathbf{x}^l$  with  $\mathbf{x}^l \in \mathbb{R}^{H^l \times W^l \times D^l}$ . Thus, we need a triplet index set  $(i^l, j^l, d^l)$  to locate any specific element in  $\mathbf{x}^l$ . The triplet  $(i^l, j^l, d^l)$  refers to one element in  $\mathbf{x}^l$ , which is in the  $d^l$ -th channel, and at spatial location  $(i^l, j^l)$  (at the  $i^l$ -th row, and  $j^l$ -th column). In actual CNN learning, the mini-batch strategy is usually used. In that case,  $\mathbf{x}^l$  becomes an order 4 tensor in  $\mathbb{R}^{H^l \times W^l \times D^l \times N}$  where  $N$  is the mini-batch size. For simplicity we assume that  $N = 1$  in this note. The results in this section, however, are easy to adopt to mini-batch versions.

In order to simplify the notations which will appear later, we follow the zero-based indexing convention, which specifies that  $0 \leq i^l < H^l$ ,  $0 \leq j^l < W^l$ , and  $0 \leq d^l < D^l$ .

In the  $l$ -th layer, a function will transform the input  $\mathbf{x}^l$  to an output  $\mathbf{y}$ , which is also the input to the next layer. Thus, we notice that  $\mathbf{y}$  and  $\mathbf{x}^{l+1}$  in fact refers to the same object, and it is very helpful to keep this point in mind. We assume the output has size  $H^{l+1} \times W^{l+1} \times D^{l+1}$ , and an element in the output is indexed by a triplet  $(i^{l+1}, j^{l+1}, d^{l+1})$ ,  $0 \leq i^{l+1} < H^{l+1}$ ,  $0 \leq j^{l+1} < W^{l+1}$ ,  $0 \leq d^{l+1} < D^{l+1}$ .

## 5 The ReLU layer

A ReLU layer does not change the size of the input, that is,  $\mathbf{x}^l$  and  $\mathbf{y}$  share the same size. In fact, the Rectified Linear Unit (hence the name ReLU) can be regarded as a truncation performed individually for every element in the input:

$$y_{i,j,d} = \max\{0, x_{i,j,d}^l\}, \quad (12)$$

with  $0 \leq i < H^l = H^{l+1}$ ,  $0 \leq j < W^l = W^{l+1}$ , and  $0 \leq d < D^l = D^{l+1}$ .

There is **no parameter** inside a ReLU layer, hence no need for parameter learning in this layer.

Based on Equation 12, it is obvious that

$$\frac{dy_{i,j,d}}{dx_{i,j,d}^l} = \llbracket x_{i,j,d}^l > 0 \rrbracket, \quad (13)$$

where  $\llbracket \cdot \rrbracket$  is the **indicator function**, being 1 if its argument is true, and 0 otherwise.

Hence, we have

$$\left[ \frac{\partial z}{\partial \mathbf{x}^l} \right]_{i,j,d} = \begin{cases} \left[ \frac{\partial z}{\partial \mathbf{y}} \right]_{i,j,d} & \text{if } x_{i,j,d}^l > 0 \\ 0 & \text{otherwise} \end{cases}. \quad (14)$$

Note that  $\mathbf{y}$  is an alias for  $\mathbf{x}^{l+1}$ .

Strictly speaking, the function  $\max(0, x)$  is not differentiable at  $x = 0$ , hence Equation 13 is a little bit **problematic in theory**. In practice, it is not an issue and ReLU is safe to use.

The purpose of ReLU is to **increase the nonlinearity** of the CNN. Since the semantic information in an image (e.g., a person and a Husky dog sitting next to each other on a bench in a garden) is obviously a highly nonlinear mapping of pixel values in the input, we want the mapping from CNN input to its output also be highly nonlinear. The ReLU function, although simple, **is a nonlinear function**, as illustrated in Figure 2.

If we treat  $x_{i,j,d}^l$  as one of the  $H^l W^l D^l$  *features* extracted by CNN layers 1 to  $l - 1$ , which can be positive, zero or negative. For example,  $x_{i,j,d}^l$  may be positive if a region inside the input image has certain patterns (like a dog's head or a cat's head or some other patterns similar to that); and  $x_{i,j,d}^l$  is negative or zero when that region does not exhibit these patterns. **The ReLU layer will set all negative values to 0, which means that  $y_{i,j,d}^l$  will be *activated* only for images possessing these patterns at that particular region.** Intuitively, this property is useful for recognizing complex patterns and objects. For example, it is only a weak evidence to support "the input image contains a cat" if a feature is activated and that feature's pattern looks like cat's head. However, if we find many activated features after the ReLU layer whose target patterns correspond to cat's head, torso, fur, legs, etc., we have higher confidence (at layer  $l + 1$ ) to say that a cat probably exists in the input image.

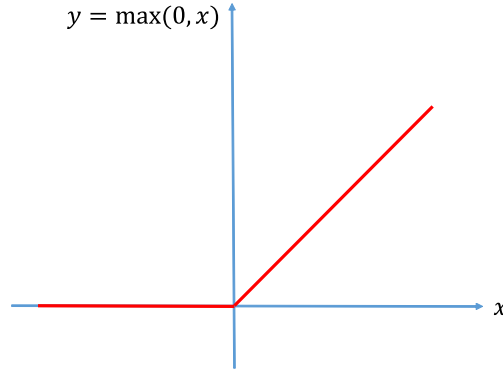


Figure 2: The ReLU function.

Other nonlinear transformations have been used in the neural network research to produce nonlinearity, for example, the logistic sigmoid function  $y = \sigma(x) = \frac{1}{1+\exp(-x)}$ . However, logistic sigmoid works significantly worse than ReLU in CNN learning. Note that  $0 < y < 1$  if a sigmoid function is used, and  $\frac{dy}{dx} = y(1 - y)$ , we have  $\frac{dy}{dx} \leq \frac{1}{4}$ . Hence, in the error back propagation process, the gradient  $\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{dy}{dx}$  will have much smaller magnitude than  $\frac{\partial z}{\partial y}$  (at most  $\frac{1}{4}$ ). In other words, a sigmoid layer will cause the magnitude of the gradient to significantly reduce, and after several sigmoid layers, the gradient will *vanish* (i.e., all its components will be close to 0). A vanishing gradient makes gradient based learning (e.g., SGD) very difficult.

On the other hand, the ReLU layer sets the gradient of some features in the  $l$ -th layer to 0, but these features are not activated (i.e., we are not interested in them). For those activated features, the gradient is back propagated without any change, which is beneficial for SGD learning. The introduction of ReLU to replace sigmoid is an important change in CNN, which significantly reduces the difficulty in learning CNN parameters and improves its accuracy. There are also more complex variants of ReLU, for example, parametric ReLU and exponential linear unit.

## 6 The convolution layer

Next, we turn to the convolution layer, which is the most involved one among those we discuss in this note.

### 6.1 What is convolution?

Let us start by convolving a matrix with one single convolution kernel. Suppose the input image is  $3 \times 4$  and the convolution kernel size is  $2 \times 2$ , as illustrated in Figure 3.

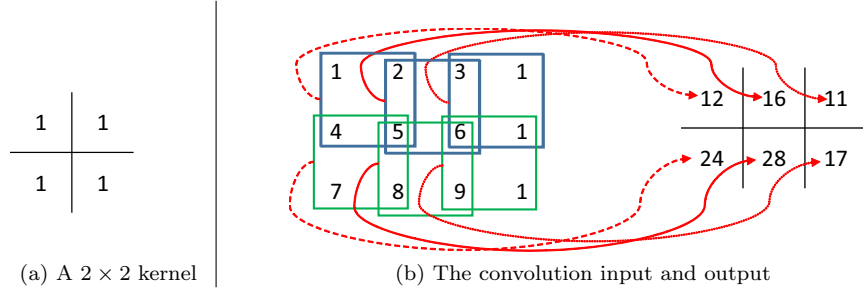


Figure 3: Illustration of the convolution operation.

If we overlap the convolution kernel on top of the input image, we can compute the product between the numbers at the same location in the kernel and the input, and we get a single number by summing these products together. For example, if we overlap the kernel with the top left region in the input, the convolution result at that spatial location is:  $1 \times 1 + 1 \times 4 + 1 \times 2 + 1 \times 5 = 12$ . We then move the kernel down by one pixel and get the next convolution result as  $1 \times 4 + 1 \times 7 + 1 \times 5 + 1 \times 8 = 24$ . We keep move the kernel down till it reaches the bottom border of the input matrix (image). Then, we return the kernel to the top, and move the kernel to its right by one element (pixel). We repeat the convolution for every possible pixel location until we have moved the kernel to the bottom right corner of the input image, as shown in Figure 3.

For order 3 tensors, the convolution operation is defined similarly. Suppose the input in the  $l$ -th layer is an order 3 tensor with size  $H^l \times W^l \times D^l$ . A convolution kernel is also an order 3 tensor with size  $H \times W \times D^l$ . When we overlap the kernel on top of the input tensor at the spatial location  $(0, 0, 0)$ , we compute the products of corresponding elements in all the  $D^l$  channels and sum the  $HW D^l$  products to get the convolution result at this spatial location. Then, we move the kernel from top to bottom and from left to right to complete the convolution.

In a convolution layer, **multiple convolution kernels** are usually used. Assuming  $D$  kernels are used and each kernel is of spatial span  $H \times W$ , we denote all the kernels as  $\mathbf{f}$ .  $\mathbf{f}$  is an order 4 tensor in  $\mathbb{R}^{H \times W \times D^l \times D}$ . Similarly, we use index variables  $0 \leq i < H$ ,  $0 \leq j < W$ ,  $0 \leq d^l < D^l$  and  $0 \leq d < D$  to pinpoint a specific element in the kernels. Also note that the set of kernels  $\mathbf{f}$  refers to the same object as the notation  $\mathbf{w}^l$  in Equation 5. We change the notation a bit to make the derivation a little bit simpler. It is also clear that even if the mini-batch strategy is used, the kernels remain unchanged.

As shown in Figure 3, the spatial extent of the output is smaller than that of the input so long as the convolution kernel is larger than  $1 \times 1$ . Sometimes we need the input and output images to have the same height and width, and a simple **padding trick** can be used. If the input is  $H^l \times W^l \times D^l$  and the kernel size is  $H \times W \times D^l \times D$ , the convolution result has size  $(H^l - H + 1) \times (W^l - W + 1) \times D$ . For every channel of the input, if we *pad* (i.e., insert)  $\lfloor \frac{H-1}{2} \rfloor$  rows above the first

row and  $\lfloor \frac{H}{2} \rfloor$  rows below the last row, and pad  $\lfloor \frac{W-1}{2} \rfloor$  columns to the left of the first column and  $\lfloor \frac{W}{2} \rfloor$  columns to the right of the last column of the input, the convolution output will be  $H^l \times W^l \times D$  in size, i.e., having the same spatial extent as the input.  $\lfloor \cdot \rfloor$  is the floor functions. Elements of the padded rows and columns are usually set to 0, but other values are also possible.

*Stride* is another important concept in convolution. In Figure 3, we convolve the kernel with the input at every possible spatial location, which corresponds to the stride  $s = 1$ . However, if  $s > 1$ , every movement of the kernel skip  $s - 1$  pixel locations (i.e., the convolution is performed once every  $s$  pixels both horizontally and vertically).

In this section, we consider the simple case when the stride is 1 and no padding is used. Hence, we have  $\mathbf{y}$  (or  $\mathbf{x}^{l+1}$ ) in  $\mathbb{R}^{H^{l+1} \times W^{l+1} \times D^{l+1}}$ , with  $H^{l+1} = H^l - H + 1$ ,  $W^{l+1} = W^l - W + 1$ , and  $D^{l+1} = D$ .

In precise mathematics, the convolution procedure can be expressed as an equation:

$$y_{i^{l+1}, j^{l+1}, d} = \sum_{i=0}^{H^l} \sum_{j=0}^{W^l} \sum_{d^l=0}^{D^l} f_{i,j,d^l,d} \times x_{i^{l+1}+i, j^{l+1}+j, d^l}^l. \quad (15)$$

Equation 15 is repeated for all  $0 \leq d \leq D = D^{l+1}$ , and for any spatial location  $(i^{l+1}, j^{l+1})$  satisfying  $0 \leq i^{l+1} < H^l - H + 1 = H^{l+1}$ ,  $0 \leq j^{l+1} < W^l - W + 1 = W^{l+1}$ . In this equation,  $x_{i^{l+1}+i, j^{l+1}+j, d^l}^l$  refers to the element of  $\mathbf{x}^l$  indexed by the triplet  $(i^{l+1} + i, j^{l+1} + j, d^l)$ .

A bias term  $b_d$  is usually added to  $y_{i^{l+1}, j^{l+1}, d}$ . We omit this term in this note for clearer presentation.

## 6.2 Why to convolve?

Figure 4 shows a color input image (4a) and its convolution results using two different kernels (4b and 4c). A  $3 \times 3$  convolution matrix  $K = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$  is used. The convolution kernel should be of size  $3 \times 3 \times 3$ , in which we set every channel to  $K$ . When there is a horizontal edge at location  $(x, y)$  (i.e., when the pixels at spatial location  $(x + 1, y)$  and  $(x - 1, y)$  differ by a large amount), we expect the convolution result to have high magnitude. As shown in Figure 4b, the convolution results indeed highlight the horizontal edges. When we set every channel of the convolution kernel to  $K^T$  (the transpose of  $K$ ), the convolution result amplifies vertical edges, as shown in Figure 4c. The matrix (or filter)  $K$  and  $K^T$  are called the Sobel operators.<sup>2</sup>

If we add a bias term to the convolution operation, we can make the convolution result positive at horizontal (vertical) edges in a certain direction (e.g., a horizontal edge with the pixels above it brighter than the pixels below it), and negative at other locations. If the next layer is a ReLU layer, the output of the next layer in fact defines many “edge detection features”, which activate

<sup>2</sup>The Sobel operator is named after Irwin Sobel, an American researcher in digital image processing.

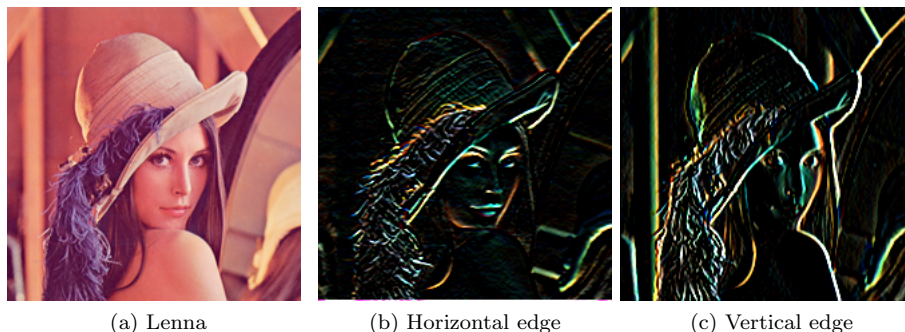


Figure 4: The Lenna image and the effect of different convolution kernels.

only at horizontal or vertical edges in certain directions. If we replace the Sobel kernel by other kernels (e.g., those learned by SGD), we can learn features that activate for **edges with different angles**. When we move further down in the deep network, subsequent layers can learn to **activate only for specific (but more complex) patterns**, e.g., groups of edges that form a particular shape. These more complex patterns will be further assembled by deeper layers to activate for semantically meaningful object parts or even a particular type of object, e.g., dog, cat, tree, beach, etc.

One more benefit of the convolution layer is that all spatial locations **share the same convolution kernel, which greatly reduces the number of parameters needed for a convolution layer**. For example, if multiple dogs appear in an input image, the same “dog-head-like pattern” feature will be activated at multiple locations, corresponding to heads of different dogs.

In a deep neural network setup, convolution also **encourages parameter sharing**. For example, suppose “dog-head-like pattern” and “cat-head-like pattern” are two features learned by a deep convolutional network. The CNN does not need to devote two sets of disjoint parameters (e.g., convolution kernels in multiple layers) for them. The CNN’s bottom layers can learn “eye-like pattern” and “animal-fur-texture pattern”, which are shared by both these more abstract features. In short, the combination of convolution kernels and deep and hierarchical structures are very effective in learning good representations (features) from images for visual recognition tasks.

We want to add a note here. Although we have used phrases such as “dog-head-like pattern”, the representation or feature learned by a CNN may not correspond exactly to semantic concepts such as “dog’s head”. A CNN feature may activate frequently for dogs’ heads and often be deactivated for other types of patterns. However, there are also possible false activations at other locations, and possible deactivations at dogs’ heads.

In fact, a key concept in CNN (or more generally deep learning) is *distributed representation*. For example, suppose our task is to recognize  $N$  different types of objects and a CNN extracts  $M$  features from any input image. It is most

likely that any one of the  $M$  features is useful for recognizing all  $N$  object categories; and to recognize one object type requires the joint effort of all  $M$  features.

### 6.3 Convolution as matrix product

Equation 15 seems pretty complex. There is a way to expand  $\mathbf{x}^l$  and simplify the convolution as a matrix product.

Let's consider a special case with  $D^l = D = 1$ ,  $H = W = 2$ , and  $H^l = 3$ ,  $W^l = 4$ . That is, we consider convolving a small single channel  $3 \times 4$  matrix (or image) with one  $2 \times 2$  filter. Using the example in Figure 3, we have

$$\begin{bmatrix} 1 & 2 & 3 & 1 \\ 4 & 5 & 6 & 1 \\ 7 & 8 & 9 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 12 & 16 & 11 \\ 24 & 28 & 17 \end{bmatrix}, \quad (16)$$

where the first matrix is denoted as  $A$ , and  $*$  is the convolution operator.

Now let's run a Matlab command `B=im2col(A,[2 2])`, we arrive at a  $B$  matrix that is an expanded version of  $A$ :

$$B = \begin{bmatrix} 1 & 4 & 2 & 5 & 3 & 6 \\ 4 & 7 & 5 & 8 & 6 & 9 \\ 2 & 5 & 3 & 6 & 1 & 1 \\ 5 & 8 & 6 & 9 & 1 & 1 \end{bmatrix}.$$

It is obvious that the first column of  $B$  corresponds to the first  $2 \times 2$  region in  $A$ , in a column-first order, corresponding to  $(i^{l+1}, j^{l+1}) = (0, 0)$ . Similarly, the second to last column in  $B$  correspond to regions in  $A$  with  $(i^{l+1}, j^{l+1})$  being  $(1, 0)$ ,  $(0, 1)$ ,  $(1, 1)$ ,  $(0, 2)$  and  $(1, 2)$ , respectively. That is, the Matlab `im2col` function explicitly expands the required elements for performing each individual convolution into a column in the matrix  $B$ . The transpose of  $B$ ,  $B^T$ , is called the `im2row` expansion of  $A$ .

Now, if we vectorize the convolution kernel itself into a vector (in the same column-first order)  $(1, 1, 1, 1)^T$ , we find that<sup>3</sup>

$$B^T \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 12 \\ 24 \\ 16 \\ 28 \\ 11 \\ 17 \end{bmatrix}. \quad (17)$$

<sup>3</sup>The notation and presentation of this note is heavily affected by the MatConvNet software package's manual (<http://arxiv.org/abs/1412.4564>, which is Matlab based). The transpose of an `im2col` expansion is equivalent to an `im2row` expansion, in which the numbers involved in one convolution is one row in the `im2row` expanded matrix. The derivation in this section uses `im2row`, complying with the implementation in MatConvNet. Caffe, a widely used CNN software package (<http://caffe.berkeleyvision.org/>, which is C++ based) uses `im2col`. These formulations are mathematically equivalent to each other.

If we reshape this resulting vector in Equation 17 properly, we get the exact convolution result matrix in Equation 16. That is, the convolution operator is a linear one. We can multiply the expanded input matrix and the vectorized filter to get a result vector, and by reshaping this vector properly we get the correct convolution results.

We can generalize this idea to more complex situations and formalize them. If  $D^l > 1$  (that is, the input  $\mathbf{x}^l$  has more than one channels), the expansion operator could first expand the first channel of  $\mathbf{x}^l$ , then the second,  $\dots$ , till all  $D^l$  channels are expanded. The expanded channels will be stacked together, that is, one row in the `im2row` expansion will have  $H \times W \times D^l$  elements, rather than  $H \times W$ .

More formally, suppose  $\mathbf{x}^l$  is a third order tensor in  $\mathbb{R}^{H^l \times W^l \times D^l}$ , with one element in  $\mathbf{x}^l$  being indexed by a triplet  $(i^l, j^l, d^l)$ . We also consider a set of convolution kernels  $\mathbf{f}$ , whose spatial extent are all  $H \times W$ . Then, the expansion operator (`im2row`) converts  $\mathbf{x}^l$  into a matrix  $\phi(\mathbf{x}^l)$ . We use two indexes  $(p, q)$  to index an element in this matrix. The expansion operator copies the element at  $(i^l, j^l, d^l)$  in  $\mathbf{x}^l$  to the  $(p, q)$ -th entry in  $\phi(\mathbf{x}^l)$ .

From the description of the expansion process, it is clear that given a fixed  $(p, q)$ , we can calculate its corresponding  $(i^l, j^l, d^l)$  triplet, because obviously

$$p = i^{l+1} + (H^l - H + 1) \times j^{l+1}, \quad (18)$$

$$q = i + H \times j + H \times W \times d^l, \quad (19)$$

$$i^l = i^{l+1} + i, \quad (20)$$

$$j^l = j^{l+1} + j. \quad (21)$$

In Equation 19, dividing  $q$  by  $HW$  and take the integer part of the quotient, we can determine which channel ( $d^l$ ) does it belong to. Similarly, we can get the offsets inside the convolution kernel as  $(i, j)$ , in which  $0 \leq i < H$  and  $0 \leq j < W$ .  $q$  completely determines one specific location inside the convolution kernel by the triplet  $(i, j, d^l)$ .

Note that the convolution result is  $\mathbf{x}^{l+1}$ , whose spatial extent is  $H^{l+1} = H^l - H + 1$  and  $W^{l+1} = W^l - W + 1$ . Thus, in Equation 18, the remainder and quotient of dividing  $p$  by  $H^{l+1} = H^l - H + 1$  will give us the offset in the convolved result  $(i^{l+1}, j^{l+1})$ , or, the top-left spatial location of the region in  $\mathbf{x}^l$  (which is to be convolved with the kernel).

Based on the definition of convolution, it is clear that we can use Equations 20 and 21 to find the offset in the input  $\mathbf{x}^l$  as  $i^l = i^{l+1} + i$  and  $j^l = j^{l+1} + j$ . That is, the mapping from  $(p, q)$  to  $(i^l, j^l, d^l)$  is one-to-one. However, we want to emphasize that the reverse mapping from  $(i^l, j^l, d^l)$  to  $(p, q)$  is one-to-many, a fact that is useful in deriving the back propagation rules in a convolution layer.

Now we use the standard `vec` operator to convert the set of convolution kernels  $\mathbf{f}$  (order 4 tensor) into a matrix. Let's start from one kernel, which can be vectorized into a vector in  $\mathbb{R}^{HWD^l}$ . Thus, all convolution kernels can be reshaped into a matrix with  $HWD^l$  rows and  $D$  columns (remember that  $D^{l+1} = D$ .) Let's call this matrix  $F$ .



Finally, with all these notations, we have a beautiful equation to calculate convolution results (cf. Equation 17, in which  $\phi(\mathbf{x}^l)$  is  $B^T$ ):

$$\text{vec}(\mathbf{y}) = \text{vec}(\mathbf{x}^{l+1}) = \text{vec}(\phi(\mathbf{x}^l)F). \quad (22)$$

Note that  $\text{vec}(\mathbf{y}) \in \mathbb{R}^{H^{l+1}W^{l+1}D}$ ,  $\phi(\mathbf{x}^l) \in \mathbb{R}^{(H^{l+1}W^{l+1}) \times (HWD^l)}$ , and  $F \in \mathbb{R}^{(HWD^l) \times D}$ . The matrix multiplication  $\phi(\mathbf{x}^l)F$  results in a matrix of size  $(H^{l+1}W^{l+1}) \times D$ . The vectorization of this resultant matrix generates a vector in  $\mathbb{R}^{H^{l+1}W^{l+1}D}$ , which matches the dimensionality of  $\text{vec}(\mathbf{y})$ .

## 6.4 The Kronecker product

A short detour to the Kronecker product is needed to compute the derivatives.

Given two matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{p \times q}$ , the Kronecker product  $A \otimes B$  is a  $mp \times nq$  matrix, defined as a block matrix

$$A \otimes B = \begin{bmatrix} a_{11}B & \cdots & a_{1n}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \cdots & a_{mn}B \end{bmatrix}. \quad (23)$$

The Kronecker product has the following properties that will be useful for us:

$$(A \otimes B)^T = A^T \otimes B^T, \quad (24)$$

$$\text{vec}(AXB) = (B^T \otimes A) \text{vec}(X), \quad (25)$$

for matrices  $A$ ,  $X$ , and  $B$  with proper dimensions (e.g., when the matrix multiplication  $AXB$  is defined.) Note that Equation 25 can be utilized from both directions.

With the help of  $\otimes$ , we can write down

$$\text{vec}(\mathbf{y}) = \text{vec}(\phi(\mathbf{x}^l)FI) = (I \otimes \phi(\mathbf{x}^l)) \text{vec}(F), \quad (26)$$

$$\text{vec}(\mathbf{y}) = \text{vec}(I\phi(\mathbf{x}^l)F) = (F^T \otimes I) \text{vec}(\phi(\mathbf{x}^l)), \quad (27)$$

where  $I$  is an identity matrix of proper size. In Equation 26, the size of  $I$  is determined the number of columns in  $F$ , hence  $I \in \mathbb{R}^{D \times D}$  in Equation 26. Similarly, in Equation 27,  $I \in \mathbb{R}^{(H^{l+1}W^{l+1}) \times (H^{l+1}W^{l+1})}$ .

The derivation for gradient computation rules in a convolution layer involves many variables and notations. We summarize the variables used in this derivation in Table 1. Note that some of these notations have not been introduced yet.

## 6.5 Backward propagation: update the parameters

As previously mentioned, we need to compute two derivatives:  $\frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)}$  and  $\frac{\partial z}{\partial \text{vec}(F)}$ , where the first term  $\frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)}$  will be used for backward propagation

Table 1: Variables, their sizes and meanings. Note that “alias” means a variable has a different name or can be reshaped into another form.

	Alias	Size & Meaning
$X$	$\mathbf{x}^l$	$H^l W^l \times D^l$ , the input tensor
$F$	$\mathbf{f}, \mathbf{w}^l$	$H W D^l \times D$ , $D$ kernels, each $H \times W$ and $D^l$ channels
$Y$	$\mathbf{y}, \mathbf{x}^{l+1}$	$H^{l+1} W^{l+1} \times D^{l+1}$ , the output, $D^{l+1} = D$
$\phi(\mathbf{x}^l)$		$H^{l+1} W^{l+1} \times H W D^l$ , the <code>im2row</code> expansion of $\mathbf{x}^l$
$M$		$H^{l+1} W^{l+1} H W D^l \times H^l W^l D^l$ , the indicator matrix for $\phi(\mathbf{x}^l)$
$\frac{\partial z}{\partial Y}$	$\frac{\partial z}{\partial \text{vec}(\mathbf{y})}$	$H^{l+1} W^{l+1} \times D^{l+1}$ , gradient for $\mathbf{y}$
$\frac{\partial z}{\partial F}$	$\frac{\partial z}{\partial \text{vec}(\mathbf{f})}$	$H W D^l \times D$ , gradient to update the convolution kernels
$\frac{\partial z}{\partial X}$	$\frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)}$	$H^l W^l \times D^l$ , gradient for $\mathbf{x}^l$ , useful for back propagation

to the previous  $((l-1)$ -th) layer, and the second term will determine how the parameters of the current  $(l)$ -th layer will be updated. A friendly reminder is to remember that  $\mathbf{f}$ ,  $F$  and  $\mathbf{w}^i$  refer to the same thing (modulo reshaping of the vector or matrix or tensor). Similarly, we can reshape  $\mathbf{y}$  into a matrix  $Y \in \mathbb{R}^{(H^{l+1} W^{l+1}) \times D}$ , then  $\mathbf{y}$ ,  $Y$  and  $\mathbf{x}^{l+1}$  refer to the same object (again modulo reshaping).

From the chain rule (Equation 10), it is easy to compute  $\frac{\partial z}{\partial \text{vec}(F)}$  as

$$\frac{\partial z}{\partial (\text{vec}(F))^T} = \frac{\partial z}{\partial (\text{vec}(Y)^T)} \frac{\partial \text{vec}(\mathbf{y})}{\partial (\text{vec}(F))^T}. \quad (28)$$

The first term in the RHS is already computed in the  $(l+1)$ -th layer as (equivalently)  $\frac{\partial z}{\partial (\text{vec}(\mathbf{x}^{l+1}))^T}$ . The second term, based on Equation 26, is pretty straightforward:

$$\frac{\partial \text{vec}(\mathbf{y})}{\partial (\text{vec}(F)^T)} = \frac{\partial ((I \otimes \phi(\mathbf{x})) \text{vec}(F))}{\partial (\text{vec}(F)^T)} = I \otimes \phi(\mathbf{x}^l). \quad (29)$$

Note that we have used the fact  $\frac{\partial X \mathbf{a}^T}{\partial \mathbf{a}} = X$  or  $\frac{\partial X \mathbf{a}}{\partial \mathbf{a}^T} = X$  so long as the matrix multiplications are well defined. This equation leads to

$$\frac{\partial z}{\partial (\text{vec}(F))^T} = \frac{\partial z}{\partial (\text{vec}(\mathbf{y})^T)} (I \otimes \phi(\mathbf{x}^l)). \quad (30)$$

Making a transpose, we get

$$\frac{\partial z}{\partial \text{vec}(F)} = (I \otimes \phi(\mathbf{x}^l))^T \frac{\partial z}{\partial \text{vec}(\mathbf{y})} \quad (31)$$

$$= (I \otimes \phi(\mathbf{x}^l)^T) \text{vec} \left( \frac{\partial z}{\partial Y} \right) \quad (32)$$

$$= \text{vec} \left( \phi(\mathbf{x}^l)^T \frac{\partial z}{\partial Y} I \right) \quad (33)$$

$$= \text{vec} \left( \phi(\mathbf{x}^l)^T \frac{\partial z}{\partial Y} \right). \quad (34)$$

Note that both Equation 25 (from RHS to LHS) and Equation 24 are used in the above derivation.

Thus, we conclude that

$$\frac{\partial z}{\partial F} = \phi(\mathbf{x}^l)^T \frac{\partial z}{\partial Y}, \quad (35)$$

which is a simple rule to update the parameters in the  $l$ -th layer: the gradient with respect to the convolution parameters is the product between  $\phi(\mathbf{x}^l)^T$  (the `im2col` expansion) and  $\frac{\partial z}{\partial Y}$  (the supervision signal transferred from the  $(l+1)$ -th layer).

## 6.6 Even higher dimensional indicator matrices

The function  $\phi(\cdot)$  has been very useful in our analysis. It is pretty high dimensional, e.g.,  $\phi(\mathbf{x}^l)$  has  $H^{l+1}W^{l+1}HWD^l$  elements. From the above, we know that an element in  $\phi(\mathbf{x}^l)$  is indexed by a pair  $p$  and  $q$ .

A quick recap about  $\phi(\mathbf{x}^l)$ : 1) from  $q$  we can determine  $d^l$ , which channel of the convolution kernel is used; and can also determine  $i$  and  $j$ , the spatial offsets inside the kernel; 2) from  $p$  we can determine  $i^{l+1}$  and  $j^{l+1}$ , the spatial offsets inside the convolved result  $\mathbf{x}^{l+1}$ ; and, 3) the spatial offsets in the input  $\mathbf{x}^l$  can be determined as  $i^l = i^{l+1} + i$  and  $j^l = j^{l+1} + j$ .

That is, the mapping  $m : (p, q) \mapsto (i^l, j^l, d^l)$  is one-to-one, and thus is a valid function. The inverse mapping, however, is one-to-many (thus not a valid function). If we use  $m^{-1}$  to represent the inverse mapping, we know that  $m^{-1}(i^l, j^l, d^l)$  is a set  $S$ , where each  $(p, q) \in S$  satisfies that  $m(p, q) = (i^l, j^l, d^l)$ .

Now we take a look at  $\phi(\mathbf{x}^l)$  from a different perspective. In order to fully specify  $\phi(\mathbf{x}^l)$ , what information is required? It is obvious that the following three types of information are needed (and only those). For *every* element of  $\phi(\mathbf{x}^l)$ , we need to know

- (A) Which region does it belong to, i.e., what is the value of  $p$  ( $0 \leq p < H^{l+1}W^{l+1}$ )?
- (B) Which element is it inside the region (or equivalently inside the convolution kernel), i.e., what is the value of  $q$  ( $0 \leq q < HWD^l$ )?

The above two types of information determines a location  $(p, q)$  inside  $\phi(\mathbf{x}^l)$ . The only missing information is

- (C) What is the value in that position, i.e.,  $[\phi(\mathbf{x}^l)]_{pq}$ ?

Since every element in  $\phi(\mathbf{x}^l)$  is a verbatim copy of one element from  $\mathbf{x}^l$ , we can turn [C] into a different but equivalent one:

- (C.1) For  $[\phi(\mathbf{x}^l)]_{pq}$ , where is this value copied from? Or, what is its original location inside  $\mathbf{x}^l$ , i.e., an index  $u$  that satisfies  $0 \leq u < H^lW^lD^l$ ?
- (C.2) The entire  $\mathbf{x}^l$ .

It is easy to see that the collective information in [A, B, C.1] (for the entire range of  $p$ ,  $q$  and  $u$ ), and [C.2] ( $\mathbf{x}^l$ ) contains exactly the same amount of information as  $\phi(\mathbf{x}^l)$ .

Since  $0 \leq p < H^{l+1}W^{l+1}$ ,  $0 \leq q < HWD^l$ , and  $0 \leq u < H^lW^lD^l$ , we can use a matrix  $M \in \mathbb{R}^{(H^{l+1}W^{l+1}HWD^l) \times (H^lW^lD^l)}$  to encode the information in [A, B, C.1]. One row index of this matrix corresponds to one location inside  $\phi(\mathbf{x}^l)$  (i.e., a  $(p, q)$  pair). One row of  $M$  has  $H^lW^lD^l$  elements, and each element can be indexed by  $(i^l, j^l, d^l)$ . Thus, each element in this matrix is indexed by a 5-tuple:  $(p, q, i^l, j^l, d^l)$ .

Then, we can use the “indicator” method to encode the function  $m(p, q) = (i^l, j^l, d^l)$  into  $M$ . That is, for any possible element in  $M$ , its row index  $x$  determines a  $(p, q)$  pair, and its column index  $y$  determines a  $(i^l, j^l, d^l)$  triplet, and  $M$  is defined as

$$M(x, y) = \begin{cases} 1 & \text{if } m(p, q) = (i^l, j^l, d^l) \\ 0 & \text{otherwise} \end{cases}. \quad (36)$$

The  $M$  matrix has the following properties:

- It is very high dimensional;
- But it is also very sparse: there is only 1 non-zero entry in the  $H^lW^lD^l$  elements in one row, because  $m$  is a function;
- $M$ , which uses information [A, B, C.1], only encodes the one-to-one correspondence between any element in  $\phi(\mathbf{x}^l)$  and any element in  $\mathbf{x}^l$ , it does not encode any specific value in  $\mathbf{x}^l$ ;
- Most importantly, putting together the one-to-one correspondence information in  $M$  and the value information in  $\mathbf{x}^l$ , obviously we have

$$\text{vec}(\phi(\mathbf{x}^l)) = M \text{vec}(\mathbf{x}^l). \quad (37)$$

## 6.7 Backward propagation: prepare supervision signal for the previous layer

In the  $l$ -th layer, we still need to compute  $\frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)}$ . For this purpose, we want to reshape  $\mathbf{x}^l$  into a matrix  $X \in \mathbb{R}^{(H^lW^l) \times D^l}$ , and use these two equivalent forms (modulo reshaping) interchangeably.

The chain rule states that  $\frac{\partial z}{\partial (\text{vec}(\mathbf{x}^l)^T)} = \frac{\partial z}{\partial (\text{vec}(\mathbf{y})^T)} \frac{\partial \text{vec}(\mathbf{y})}{\partial (\text{vec}(\mathbf{x}^l)^T)}$  (cf. Equation 11). We will start by studying the second term in the RHS (utilizing Equations 27 and 37):

$$\frac{\partial \text{vec}(\mathbf{y})}{\partial (\text{vec}(\mathbf{x}^l)^T)} = \frac{\partial (F^T \otimes I) \text{vec}(\phi(\mathbf{x}^l))}{\partial (\text{vec}(\mathbf{x}^l)^T)} = (F^T \otimes I)M. \quad (38)$$

Thus,

$$\frac{\partial z}{\partial (\text{vec}(\mathbf{x}^l)^T)} = \frac{\partial z}{\partial (\text{vec}(\mathbf{y})^T)} (F^T \otimes I)M. \quad (39)$$

Since (using Equation 25 from right to left)

$$\frac{\partial z}{\partial(\text{vec}(\mathbf{y})^T)}(F^T \otimes I) = \left( (F \otimes I) \frac{\partial z}{\partial \text{vec}(\mathbf{y})} \right)^T \quad (40)$$

$$= \left( (F \otimes I) \text{vec} \left( \frac{\partial z}{\partial Y} \right) \right)^T \quad (41)$$

$$= \text{vec} \left( I \frac{\partial z}{\partial Y} F^T \right)^T \quad (42)$$

$$= \text{vec} \left( \frac{\partial z}{\partial Y} F^T \right)^T, \quad (43)$$

we have

$$\frac{\partial z}{\partial(\text{vec}(\mathbf{x}^l)^T)} = \text{vec} \left( \frac{\partial z}{\partial Y} F^T \right)^T M, \quad (44)$$

or equivalently

$$\frac{\partial z}{\partial(\text{vec}(\mathbf{x}^l))} = M^T \text{vec} \left( \frac{\partial z}{\partial Y} F^T \right). \quad (45)$$

Let's have a closer look at the RHS.  $\frac{\partial z}{\partial Y} F^T \in \mathbb{R}^{(H^{l+1}W^{l+1}) \times (HWD^l)}$ , and  $\text{vec} \left( \frac{\partial z}{\partial Y} F^T \right)$  is a vector in  $\mathbb{R}^{H^{l+1}W^{l+1}HWD^l}$ . On the other hand,  $M^T$  is an indicator matrix in  $\mathbb{R}^{(H^lW^lD^l) \times (H^{l+1}W^{l+1}HWD^l)}$ .

In order to pinpoint one element in  $\text{vec}(\mathbf{x}^l)$  or one row in  $M^T$ , we need an index triplet  $(i^l, j^l, d^l)$ , with  $0 \leq i^l < H^l$ ,  $0 \leq j^l < W^l$ , and  $0 \leq d^l < D^l$ . Similarly, to locate a column in  $M^T$  or an element in  $\frac{\partial z}{\partial Y} F^T$ , we need an index pair  $(p, q)$ , with  $0 \leq p < H^{l+1}W^{l+1}$  and  $0 \leq q < HWD^l$ .

Thus, the  $(i^l, j^l, d^l)$ -th entry of  $\frac{\partial z}{\partial(\text{vec}(\mathbf{x}^l))}$  equals the multiplication of two vectors: the row in  $M^T$  (or the column in  $M$ ) that is indexed by  $(i^l, j^l, d^l)$ , and  $\text{vec} \left( \frac{\partial z}{\partial Y} F^T \right)$ .

Furthermore, since  $M^T$  is an indicator matrix, in the row vector indexed by  $(i^l, j^l, d^l)$ , only those entries whose index  $(p, q)$  satisfies  $m(p, q) = (i^l, j^l, d^l)$  have a value 1, all other entries are 0. Thus, the  $(i^l, j^l, d^l)$ -th entry of  $\frac{\partial z}{\partial(\text{vec}(\mathbf{x}^l))}$  equals the sum of these corresponding entries in  $\text{vec} \left( \frac{\partial z}{\partial Y} F^T \right)$ .

Transferring the above description into precise mathematical form, we get the following succinct equation:

$$\left[ \frac{\partial z}{\partial X} \right]_{(i^l, j^l, d^l)} = \sum_{(p, q) \in m^{-1}(i^l, j^l, d^l)} \left[ \frac{\partial z}{\partial Y} F^T \right]_{(p, q)}. \quad (46)$$

In other words, to compute  $\frac{\partial z}{\partial X}$ , we do not need to explicitly use the extremely high dimensional matrix  $M$ . Instead, Equation 46 and Equations 18 to 21 can be used to efficiently find  $\frac{\partial z}{\partial X}$ .

We use the simple convolution example in Figure 3 to illustrate the inverse mapping  $m^{-1}$ , which is shown in Figure 5.

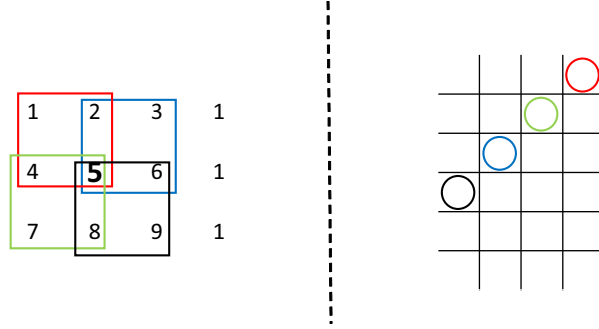


Figure 5: Illustration of how to compute  $\frac{\partial z}{\partial X}$ .

In the right half of Figure 5, the  $6 \times 4$  matrix is  $\frac{\partial z}{\partial Y} F^T$ . In order to compute the partial derivative of  $z$  with respect to one element in the input  $X$ , we need to find which elements in  $\frac{\partial z}{\partial Y} F^T$  is involved and add them. In the left half of Figure 5, we show that the input element 5 (shown in larger font) is involved in 4 convolution operations, shown by the red, green, blue and black boxes, respectively. These 4 convolution operations correspond to  $p = 1, 2, 3, 4$ . For example, when  $p = 2$  (the green box), 5 is the third element in the convolution, hence  $q = 3$  when  $p = 2$  and we put a green circle in the  $(2, 3)$ -th element of the  $\frac{\partial z}{\partial Y} F^T$  matrix. After all 4 circles are put in the  $\frac{\partial z}{\partial Y} F^T$  matrix, the partial derivative is the sum of elements in these four locations of  $\frac{\partial z}{\partial Y} F^T$ .

The set  $m^{-1}(i^l, j^l, d^l)$  contains at most  $HW D^l$  elements. Hence, Equation 46 requires at most  $HW D^l$  summations to compute one element of  $\frac{\partial z}{\partial X}$ .<sup>4</sup>

## 6.8 Fully connected layer as a convolution layer

As aforementioned, one benefit of the convolution layer is that convolution is a local operation. The spatial extent of a kernel is often small (e.g.,  $3 \times 3$ ). One element in  $\mathbf{x}^{l+1}$  is usually computed using only a small number of elements in its input  $\mathbf{x}^l$ .

A fully connected layer refers to a layer if the computation of any element in the output  $\mathbf{x}^{l+1}$  (or  $\mathbf{y}$ ) requires all elements in the input  $\mathbf{x}^l$ . A fully connected layer is sometimes useful at the end of a deep CNN model. For example, if after many convolution, ReLU and pooling (which will be discussed soon) layers, the output of the current layer contain distributed representations for the input image, we want to use all these features in the current layer to build features with stronger capabilities in the next one. A fully connected layer is useful for this purpose.

Suppose the input of a layer  $\mathbf{x}^l$  has size  $H^l \times W^l \times D^l$ . If we use convolution kernels whose size is  $H^l \times W^l \times D^l$ , then  $D$  such kernels form an order 4 tensor

<sup>4</sup>In Caffe, this computation is implemented by a function called `col2im`. In MatConvNet, this operation is operated in a `row2im` manner, although the name `row2im` is not explicitly used.

in  $H^l \times W^l \times D^l \times D$ . The output is  $\mathbf{y} \in \mathbb{R}^D$ . It is obvious that to compute any element in  $\mathbf{y}$ , we need to use all elements in the input  $\mathbf{x}^l$ . Hence, this layer is a fully connected layer, but can be implemented as a convolution layer. Hence, we do not need to derive learning rules for a fully connected layer separately.

## 7 The pooling layer

We will use the same notation inherited from the convolution layer. Let  $\mathbf{x}^l \in \mathbb{R}^{H^l \times W^l \times D^l}$  be the input to the  $l$ -th layer, which is now a pooling layer. The pooling operation requires no parameter (i.e.,  $\mathbf{w}^i$  is null, hence parameter learning is not needed for this layer). The spatial extent of the pooling ( $H \times W$ ) is specified in the design of the CNN structure. Assume that  $H$  divides  $H^l$  and  $W$  divides  $W^l$  and the stride equals the pooling spatial extent,<sup>5</sup> the output of pooling ( $\mathbf{y}$  or equivalently  $\mathbf{x}^{l+1}$ ) will be an order 3 tensor of size  $H^{l+1} \times W^{l+1} \times D^{l+1}$ , with

$$H^{l+1} = \frac{H^l}{H}, \quad W^{l+1} = \frac{W^l}{W}, \quad D^{l+1} = D^l. \quad (47)$$

A pooling layer operates upon  $\mathbf{x}^l$  channel by channel independently. Within each channel, the matrix with  $H^l \times W^l$  elements are divided into  $H^{l+1} \times W^{l+1}$  nonoverlapping subregions, each subregion being  $H \times W$  in size. The pooling operator then maps a subregion into a single number.

Two types of pooling operators are widely used: **max pooling** and **average pooling**. In max pooling, the pooling operator maps a subregion to its maximum value, while the average pooling maps a subregion to its average value. In precise mathematics,

$$\text{max :} \quad y_{i^{l+1}, j^{l+1}, d} = \max_{0 \leq i < H, 0 \leq j < W} x_{i^{l+1} \times H + i, j^{l+1} \times W + j, d}^l, \quad (48)$$

$$\text{average :} \quad y_{i^{l+1}, j^{l+1}, d} = \frac{1}{HW} \sum_{0 \leq i < H, 0 \leq j < W} x_{i^{l+1} \times H + i, j^{l+1} \times W + j, d}^l, \quad (49)$$

where  $0 \leq i^{l+1} < H^{l+1}$ ,  $0 \leq j^{l+1} < W^{l+1}$ , and  $0 \leq d < D^{l+1} = D^l$ .

Pooling is a local operator, and its forward computation is pretty straightforward. Now we focus on the back propagation. Only max pooling is discussed and we can resort to the indicator matrix again.<sup>6</sup> All we need to encode in this indicator matrix is: for every element in  $\mathbf{y}$ , where does it come from in  $\mathbf{x}^l$ ?

We need a triplet  $(i^l, j^l, d^l)$  to pinpoint one element in the input  $\mathbf{x}^l$ , and another triplet  $(i^{l+1}, j^{l+1}, d^{l+1})$  to locate one element in  $\mathbf{y}$ . The pooling output  $y_{i^{l+1}, j^{l+1}, d^{l+1}}$  comes from  $x_{i^l, j^l, d^l}^l$ , if and only if the following conditions are met:

- They are in the same channel;
- The  $(i^l, j^l)$ -th spatial entry belongs to the  $(i^{l+1}, j^{l+1})$ -th subregion;

<sup>5</sup>That is, the strides in the vertical and horizontal direction are  $H$  and  $W$ , respectively. The most widely used pooling setup is  $H = W = 2$  with a stride 2.

<sup>6</sup>Average pooling can be dealt with using a similar idea.

- The  $(i^l, j^l)$ -th spatial entry is the largest one in that subregion.

Translating these conditions into equations, we get

$$d^{l+1} = d^l, \quad (50)$$

$$\left\lfloor \frac{i^l}{H} \right\rfloor = i^{l+1}, \left\lfloor \frac{j^l}{W} \right\rfloor = j^{l+1}, \quad (51)$$

$$x_{i^l, j^l, d^l}^l \geq y_{i+i^{l+1} \times H, j+j^{l+1} \times W, d^l}, \forall 0 \leq i < H, 0 \leq j < W, \quad (52)$$

where  $\lfloor \cdot \rfloor$  is the floor function. If the stride is not  $H$  ( $W$ ) in the vertical (horizontal) direction, Equation 51 must be changed accordingly.

Given a  $(i^{l+1}, j^{l+1}, d^{l+1})$  triplet, there is only one  $(i^l, j^l, d^l)$  triplet that satisfies all these conditions. Thus, we define an indicator matrix

$$S(\mathbf{x}^l) \in \mathbb{R}^{(H^{l+1}W^{l+1}D^{l+1}) \times (H^lW^lD^l)}. \quad (53)$$

One triplet of indexes  $(i^{l+1}, j^{l+1}, d^{l+1})$  specifies a row in  $S$ , while  $(i^l, j^l, d^l)$  specifies a column. These two triplets together pinpoint one element in  $S(\mathbf{x}^l)$ . We set that element to 1 if Equations 50 to 52 are simultaneously satisfied, and 0 otherwise. One row of  $S(\mathbf{x}^l)$  corresponds to one element in  $\mathbf{y}$ , and one column corresponds to one element in  $\mathbf{x}^l$ .

With the help of this indicator matrix, we have

$$\text{vec}(\mathbf{y}) = S(\mathbf{x}^l) \text{vec}(\mathbf{x}^l). \quad (54)$$

Then, it is obvious that

$$\frac{\partial \text{vec}(\mathbf{y})}{\partial (\text{vec}(\mathbf{x}^l)^T)} = S(\mathbf{x}^l), \quad \frac{\partial z}{\partial (\text{vec}(\mathbf{x}^l)^T)} = \frac{\partial z}{\partial (\text{vec}(\mathbf{y})^T)} S(\mathbf{x}^l), \quad (55)$$

and consequently

$$\frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)} = S(\mathbf{x}^l)^T \frac{\partial z}{\partial \text{vec}(\mathbf{y})}. \quad (56)$$

$S(\mathbf{x}^l)$  is very sparse. It has exactly one nonzero entry in every row. Thus, we do not need to use the entire matrix in the computation. Instead, we just need to record the locations of those nonzero entries—there are only  $H^{l+1}W^{l+1}D^{l+1}$  such entries in  $S(\mathbf{x}^l)$ .

A simple example can explain the meaning of these equations. Let us consider a  $2 \times 2$  max pooling with stride 2. For a given channel  $d^l$ , the first spatial subregion contains four elements in the input, with  $(i, j) = (0, 0), (1, 0), (0, 1)$  and  $(1, 1)$ , and let us suppose the element at spatial location  $(0, 1)$  is the largest among them. In the forward pass, the value indexed by  $(0, 1, d^l)$  in the input (i.e.,  $x_{0,1,d^l}^l$ ) will be assigned to the element in the  $(0, 0, d^l)$ -th element in the output (i.e.,  $y_{0,0,d^l}$ ).

One column in  $S(\mathbf{x}^l)$  contains at most one nonzero element if the strides are  $H$  and  $W$ , respectively. In the above example, the column of  $S(\mathbf{x}^l)$  indexed by



$(0, 0, d^l)$ ,  $(1, 0, d^l)$  and  $(1, 1, d^l)$  are all zero vectors. The column corresponding to  $(0, 1, d^l)$  contains only one nonzero entry, whose row index is determined by  $(0, 0, d^l)$ . Hence, in the back propagation, we have

$$\left[ \frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)} \right]_{(0,1,d^l)} = \left[ \frac{\partial z}{\partial \text{vec}(\mathbf{y})} \right]_{(0,0,d^l)},$$

and

$$\left[ \frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)} \right]_{(0,0,d^l)} = \left[ \frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)} \right]_{(1,0,d^l)} = \left[ \frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)} \right]_{(1,1,d^l)} = 0.$$

However, if the pooling strides are smaller than  $H$  and  $W$  in the vertical and horizontal directions, respectively, one element in the input tensor may be the largest element in several pooling subregions. Hence, there can have more than one nonzero entries in one column of  $S(\mathbf{x}^l)$ . Let us consider the example input in Figure 5. If a  $2 \times 2$  max pooling is applied to it and the stride is 1 in both directions, the element 9 is the largest in two pooling regions:  $\begin{bmatrix} 5 & 6 \\ 8 & 9 \end{bmatrix}$  and  $\begin{bmatrix} 6 & 1 \\ 9 & 1 \end{bmatrix}$ . Hence, in the column of  $S(\mathbf{x}^l)$  corresponding to the element 9 (indexed by  $(2, 2, d^l)$  in the input tensor), there are two nonzero entries whose row indexes correspond to  $(i^{l+1}, j^{l+1}, d^{l+1}) = (1, 1, d^l)$  and  $(1, 2, d^l)$ . Thus, in this example, we have

$$\left[ \frac{\partial z}{\partial \text{vec}(\mathbf{x}^l)} \right]_{(2,2,d^l)} = \left[ \frac{\partial z}{\partial \text{vec}(\mathbf{y})} \right]_{(1,1,d^l)} + \left[ \frac{\partial z}{\partial \text{vec}(\mathbf{y})} \right]_{(1,2,d^l)}.$$

## 8 A case study: the VGG-16 net

We have introduced the convolution, pooling, ReLU and fully connected layers till now, and have briefly mentioned the softmax layer. With these layers, we can build many powerful deep CNN models.

### 8.1 VGG-Verydeep-16

The VGG-Verydeep-16 CNN model is a pretrained CNN model released by the Oxford VGG group.<sup>7</sup> We use it as an example to study the detailed structure of CNN networks. The VGG-16 model architecture is listed in Table 2.

There are six types of layers in this model.

**Convolution** A convolution layer is abbreviated as “Conv”. Its description includes three parts: number of channels; kernel spatial extent (kernel size); padding (‘p’) and stride (‘st’) size.

**ReLU** No description is needed for a ReLU layer.

<sup>7</sup>[http://www.robots.ox.ac.uk/~vgg/research/very\\_deep/](http://www.robots.ox.ac.uk/~vgg/research/very_deep/)

Table 2: The VGG-Verydeep-16 architecture and receptive field

	type	description	r. size		type	description	r. size
1	Conv	64;3x3;p=1,st=1	212	20	Conv	512;3x3;p=1,st=1	20
2	ReLU		210	21	ReLU		18
3	Conv	64;3x3;p=1,st=1	210	22	Conv	512;3x3;p=1,st=1	18
4	ReLU		208	23	ReLU		16
5	Pool	2x2;st=2	208	24	Pool	2x2;st=2	16
6	Conv	128;3x3;p=1,st=1	104	25	Conv	512;3x3;p=1,st=1	8
7	ReLU		102	26	ReLU		6
8	Conv	128;3x3;p=1,st=1	102	27	Conv	512;3x3;p=1,st=1	6
9	ReLU		100	28	ReLU		4
10	Pool	2x2;st=2	100	29	Conv	512;3x3;p=1,st=1	4
11	Conv	256;3x3;p=1,st=1	50	30	ReLU		2
12	ReLU		48	31	Pool		2
13	Conv	256;3x3;p=1,st=1	48	32	FC	(7x7x512)x4096	1
14	ReLU		46	33	ReLU		
15	Conv	256;3x3;p=1,st=1	46	34	Drop	0.5	
16	ReLU		44	35	FC	4096x4096	
17	Pool	2x2;st=2	44	36	ReLU		
18	Conv	512;3x3;p=1,st=1	22	37	Drop	0.5	
19	ReLU		20	38	FC	4096x1000	
				39	$\sigma$	(softmax layer)	

**Pool** A pooling layer is abbreviated as “Pool”. Only max pooling is used in VGG-16. The pooling kernel size is always  $2 \times 2$  and the stride is always 2 in VGG-16.

**Fully connected** A fully connected layer is abbreviated as “FC”. Fully connected layers are implemented using convolution in VGG-16. Its size is shown in the format  $n_1 \times n_2$ , where  $n_1$  is the size of the input tensor, and  $n_2$  is the size of the output tensor. Although  $n_1$  can be a triplet (such as  $7 \times 7 \times 512$ ,  $n_2$  is always an integer.

**Dropout** A dropout layer is abbreviated as “Drop”. Dropout is a technique to improve the generalization of deep learning methods. It sets the weights connected to a certain percentage of nodes in the network to 0 (and VGG-16 set the percentage to 0.5 in the two dropout layers).

**Softmax** It is abbreviated as “ $\sigma$ ”.

We want to add a few notes about this example deep CNN architecture.

- A convolution layer is always followed by a ReLU layer in VGG-16. The ReLU layers increase the nonlinearity of the CNN model.
- The convolution layers between two pooling layers have the same number of channels, kernel size and stride. In fact, stacking two  $3 \times 3$  convolution layers is equivalent to one  $5 \times 5$  convolution layer; and stacking three  $3 \times 3$  convolution kernels replaces a  $7 \times 7$  convolution layer. Stacking a few (2 or 3) smaller convolution kernels, however, computes faster than a large

convolution kernel. In addition, the number of parameters is also reduced, e.g.,  $2 \times 3 \times 3 = 18 < 25 = 5 \times 5$ . The ReLU layers inserted in between small convolution layers are also helpful.

- The input to VGG-16 is an image with size  $224 \times 224 \times 3$ . Because the padding is one in the convolution kernels (meaning one row or column is added outside of the four edges of the input), convolution will not change the spatial extent. The pooling layers will reduce the input size by a factor of 2. Hence, the output after the last (5th) pooling layer has spatial extent  $7 \times 7$  (and 512 channels). We may interpret this tensor as  $7 \times 7 \times 512 = 25088$  “features”. The first fully connected layer converts them into 4096 features. The number of features remains at 4096 after the second fully connected layer.
- The VGG-16 is trained for the ImageNet classification challenge, which is an object recognition problem with 1000 classes. The last fully connected layer ( $4096 \times 1000$ ) output a length 1000 vector for every input image, and the softmax layer converts this length 1000 vector into the estimated posterior probability for the 1000 classes.

## 8.2 Receptive field

Another important concept in CNN is the **receptive field size** (abbreviated as “r. size” in Table 2). Let us look at one element in the input to the first fully connected layer (32|FC). Because it is the output of a max pooling, we need values in a  $2 \times 2$  spatial extent in the input to the max pool layer to compute this element (and we only need elements in this spatial extent). This  $2 \times 2$  spatial extent is called the *receptive field* for this element. In Table 2, we **listed the spatial extent for any element in the output of the last pooling layer**. Note that because the receptive field is square, we only use one number (e.g., 48 for  $48 \times 48$ ). **The receptive field size listed for one layer is the spatial extent in the input to that layer.**

A  $3 \times 3$  convolution layer will increase the receptive field by 2 and a pooling layer will double the spatial extent. As shown in Table 2, receptive field size in the input to the first layer is  $212 \times 212$ . In other words, in order to compute any single element in the  $7 \times 7 \times 512$  output of the last pooling layer, a  $212 \times 212$  image patch is required (including the padded pixels in all convolution layers).

It is obvious that **the receptive field size increases when the network becomes deeper**, especially when a pooling layer is added to the deep net. Unlike traditional computer vision and image processing features which depend only on a small receptive field (e.g.,  $16 \times 16$ ), deep CNN computes its representation (or features) using large receptive fields. The larger receptive field characteristic is an important reason why CNN has achieved **higher accuracy** than classic methods in image recognition.

## 9 Remarks

We hope this introductory note on CNN is clear, self-contained, and easy to understand to our readers.

Once a reader is confident in his/her understanding of CNN at the mathematical level, in the next step it is very helpful to get some hands on CNN experience. For example, one can validate what has been talked about in this note using the MatConvNet software package if you prefer the Matlab environment.<sup>8</sup> For C++ lovers, Caffe is a widely used tool.<sup>9</sup> The Theano package is a python package for deep learning.<sup>10</sup> Many more resources for deep learning (not only CNN) are available, e.g., Torch,<sup>11</sup> TensorFlow,<sup>12</sup> and more will emerge soon.

## Exercises

1. Dropout is a very useful technique in training neural networks, which is proposed by Srivastava et al. in a paper titled “Dropout: A Simple Way to Prevent Neural Networks from Overfitting” in JMLR.<sup>13</sup> *Carefully* read this paper and answer the following questions (please organize your answer to every question in one brief sentence).
  - (a) How does dropout operate during training?
  - (b) How does dropout operate during testing?
  - (c) What is the benefit of dropout?
  - (d) Why dropout can achieve this benefit?
2. The VGG16 CNN model (also called VGG-Verydeep-16) was publicized by Karen Simonyan and Andrew Zisserman in a paper titled “Very Deep Convolutional Networks for Large-Scale Image Recognition” in the arXiv preprint server.<sup>14</sup> And, the GoogLeNet model was publicized by Szegedy et al. in a paper titled “Going Deeper with Convolutions” in the arXiv preprint server.<sup>15</sup> These two papers were publicized around the same time and share some similar ideas. *Carefully* read both papers and answer the following questions (please organize your answer to every question in one brief sentence).
  - (a) Why do they use small convolution kernels (mainly  $3 \times 3$ ) rather than

---

<sup>8</sup><http://www.vlfeat.org/matconvnet/>

<sup>9</sup><http://caffe.berkeleyvision.org/>

<sup>10</sup><http://deeplearning.net/software/theano/>

<sup>11</sup><http://torch.ch/>

<sup>12</sup><https://www.tensorflow.org/>

<sup>13</sup>Available at <http://jmlr.org/papers/v15/srivastava14a.html>

<sup>14</sup>Available at <https://arxiv.org/abs/1409.1556>, later published in ICLR 2015 as a conference track paper.

<sup>15</sup>Available at <https://arxiv.org/abs/1409.4842>, later published in CVPR 2015.

larger ones?

(b) Why both networks are quite deep (i.e., with many layers, around 20)?

(c) Which difficulty is caused by the large depth? How are they solved in these two networks?

3. Batch Normalization (BN) is another very useful technique in training deep neural networks, which is proposed by Sergey Ioffe and Christian Szegedy, in a paper titled “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift” in ICML 2015.<sup>16</sup> *Carefully* read this paper and answer the following questions (please organize your answer to every question in one brief sentence).

(a) What is internal covariate shift?

(b) How does BN deal with this?

(c) How does BN operate in a convolution layer?

(d) What is the benefit of using BN?

4. ResNet is a very deep neural network learning technique proposed by He et al. in a paper titled “Deep Residual Learning for Image Recognition” in CVPR 2016.<sup>17</sup> *Carefully* read this paper and answer the following questions (please organize your answer to every question in one brief sentence).

(a) Although VGG16 and GoogLeNet have encountered difficulties in training networks around 20–30 layers, what enables ResNet to train networks as deep as 1000 layers?

(b) VGG16 is a feed-forward network, where each layer has only one input and only one output. While GoogLeNet and ResNet are DAGs (directed **acyclic** graph), where **one layer can have multiple inputs and multiple outputs**, so long as the data flow in the network structure does *not* form a cycle. What is the benefit of DAG vs. feed-forward?

(c) VGG16 has two fully connected layers (fc6 and fc7), while ResNet and GoogLeNet do not have fully connected layers (except the last layer for classification). What is used to **replace FC** in them? What is the benefit?

5. *AlexNet* refers to the deep convolutional neural network trained on the ILSVRC challenge data, which is a groundbreaking work of deep CNN for computer vision tasks. The technical details of AlexNet is reported in the paper “ImageNet Classification with Deep Convolutional Neural Networks”, by Alex Krizhevsky, Ilya Sutskever and Geoffrey E. Hinton in NIPS 25.<sup>18</sup> It proposed the ReLU activation function and creatively used GPUs to accelerate the computations. *Carefully* read this paper

---

<sup>16</sup> Available at <http://jmlr.org/proceedings/papers/v37/ioffe15.pdf>

<sup>17</sup> Available at <https://arxiv.org/pdf/1512.03385.pdf>

<sup>18</sup> This paper is available at <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks>

and answer the following questions (please organize your answer to every question in one brief sentence).

- (a) Describe your understanding of how ReLU helps its success? And, how do the GPUs help out?
- (b) Using the average of predictions from several networks help reduce the error rates. Why?
- (c) Where is the dropout technique applied? How does it help? And what is the cost of using dropout?
- (d) How many parameters are there in AlexNet? Why the dataset size (1.2 million) is important for the success of AlexNet?

6. We will try different CNN structures on the MNIST dataset. We denote the “baseline” network in the MNIST example in MatConvNet as BASE in this question.<sup>19</sup> In this question, a convolution layer is denoted as “ $x \times y \times \text{nIn} \times \text{nOut}$ ”, whose kernel size is  $x \times y$ , with nIn input and nOut output channels, with stride equal 1 and pad equal 0. The pooling layers are  $2 \times 2$  max pooling with stride equal 2. The BASE network has four blocks. The first consists of a  $5 \times 5 \times 1 \times 20$  convolution and a max pooling; the second block is composed of a  $5 \times 5 \times 20 \times 50$  convolution and a max pooling; the third block is a  $4 \times 4 \times 50 \times 500$  convolution (FC) plus a ReLU layer; and the final block is the classification layer ( $1 \times 1 \times 500 \times 10$  convolution).

(a) The MNIST dataset is available at [yann.lecun.com/exdb/mnist](http://yann.lecun.com/exdb/mnist). Read the instructions in that page, and write a program to transform the data to formats that suit your favorite deep learning software.

(b) Learning deep learning models often involve random numbers. Before the training starts, set the random number generator’s seed to 0. Then, use the BASE network structure and the *first 10000 training examples* to learn its parameters. What is test set error rate (on the 10000 test examples) after 20 training epochs?

(c) From now on, if not otherwise specified, we assume the first 10000 training examples and 20 epochs are used. Now we define the BN network structure, which adds a batch normalization layer after every convolution layer in the first three blocks. What is its error rate? What will you say about BN vs. BASE?

(d) If you add a dropout layer after the classification layer in the 4th block. What is the new error rate of BASE and BN? What you will comment on dropout?

(e) Now we define the SK network structure, which refers to small kernel size. SK is based on BN. The first block ( $5 \times 5$  convolution plus pooling) now is changed to two  $3 \times 3$  convolutions, and BN + ReLU is applied after

---

<sup>19</sup>MatConvNet version 1.0-beta20. Please refer to MatConvNet for all the details of BASE, such as parameter initialization and learning rate.

every convolution. For example, block 1 is now  $3 \times 3 \times 1 \times 20$  convolution + BN + ReLU +  $3 \times 3 \times 20 \times 20$  convolution + BN + ReLU + pool. What is SK's error rate? How will you comment on that (e.g., how and why the error rate changes?)

(f) Now we define the SK-s networks structure. The notation 's' refers to a multiplier that changes the number of channels in convolution layers. For example, SK is the same as SK-1. And, SK-2 means the number of channels in all convolution layers (except the one in block 4) are multiplied by 2. Train networks for SK-2, SK-1.5, SK-1, SK-0.5 and SK-0.2. Report their error rates and comment on them.

(g) Now we experiment with different training set sizes using the SK-0.2 network structure. Using the first 500, 1000, 2000, 5000, 10000, 20000, and 60000 (all) training examples, what error rates do you achieve? Comment on your observations.

(h) Using the SK-0.2 network structure, study how different training sets affect its performance. Train 6 networks, and use the  $(10000 \times (i - 1) + 1)$ -th to  $(i \times 10000)$ -th training examples in training the  $i$ -th network. Are CNNs stable in terms of different training sets?

(i) Now we study how randomness affects CNN learning. Instead of set the random number generator's seed to 0, use 1, 12, 123, 1234, 12345 and 123456 as the seed to train 6 different SK-0.2 networks. What are their error rates? Comment on your observations.

(j) Finally, in SK-0.2, change all ReLU layers to sigmoid layers. How do you comment on the comparison on error rates of using ReLU and sigmoid activation functions?